

Worker/Wrapper/Makes it/Faster

Jennifer Hackett Graham Hutton

School of Computer Science, University of Nottingham

{jph,gmh}@cs.nott.ac.uk

Abstract

Much research in program optimization has focused on formal approaches to correctness: proving that the meaning of programs is preserved by the optimisation. Paradoxically, there has been comparatively little work on formal approaches to efficiency: proving that the performance of optimized programs is actually improved. This paper addresses this problem for a general-purpose optimization technique, the worker/wrapper transformation. In particular, we use the call-by-need variant of improvement theory to establish conditions under which the worker/wrapper transformation is formally guaranteed to preserve or improve the time performance of programs in lazy languages such as Haskell.

Categories and Subject Descriptors D.1.1 [*Programming Techniques*]: Applicative (Functional) Programming

Keywords general recursion; improvement

1. Introduction

To misquote Oscar Wilde [31], “functional programmers know the value of everything and the cost of nothing”¹. More precisely, the functional approach to programming emphasises what programs *mean* in a denotational sense, rather than what programs *do* in terms of their operational behaviour. For many programming tasks this emphasis is entirely appropriate, allowing the programmer to focus on the high-level description of what is being computed rather than the low-level details of how this is realised. However, in the context of program optimisation both aspects play a central role, as the aim of optimisation is to improve the operational performance of programs while maintaining their denotational correctness.

A research paper on program optimisation therefore should justify both the correctness and performance aspects of the optimisation described. There is a whole spectrum of possible approaches to this, ranging from informal tests and benchmarks [19], to tool-based methods such as property-

based testing [3] and space/time profiling [24], all the way up to formal mathematical proofs [17]. For correctness, it is now becoming standard to formally prove that an optimisation preserves the meaning of programs. For performance, however, the standard approach is to provide some form of empirical evidence that an optimisation improves the efficiency of programs, and there is little published work on formal proofs of improvement.

In this paper, we aim to go some way toward redressing this imbalance in the context of the *worker/wrapper transformation* [7], putting the denotational and operational aspects on an equally formal footing. The worker/wrapper transformation is a general purpose optimisation technique that has already been formally proved correct, as well as being realised in practice as an extension to the Glasgow Haskell Compiler [26]. In this paper we formally prove that this transformation is guaranteed to preserve or improve time performance with respect to an established operational theory. In other words, we show that the worker/wrapper transformation never makes programs slower. Specifically, the paper makes the following contributions:

- We show how Moran and Sands’ work on *call-by-need improvement theory* [15] can be applied to formally justify that the worker/wrapper transformation for least fixed points preserves or improves time performance;
- We present preconditions that ensure the transformation improves performance in this manner, which come naturally from the preconditions that ensure correctness;
- We demonstrate the utility of the new theory by verifying that examples from previous worker/wrapper papers indeed exhibit a time improvement.

The use of *call-by-need* improvement theory means that our work applies to lazy functional languages such as Haskell. Traditionally, the operational behaviour of lazy evaluation has been seen as difficult to reason about, but we show that with the right tools this need not be the case. To the best of our knowledge, this paper is the first time that a general purpose optimisation method for lazy languages has been formally proved to improve time performance.

Improvement theory does not seem to have attracted much attention in recent years, but we hope that this paper can help to generate more interest in this and other techniques for reasoning about lazy evaluation. Whereas in many papers calculations and proofs are often omitted or compressed for reasons of brevity, in this paper they are the central focus, so are presented in detail.

¹The general form of this misquote is due to Alan Perlis, who originally said it of Lisp programmers.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICFP ’14, September 1–6, 2014, Gothenburg, Sweden.

Copyright is held by the owner/author(s). Publication rights licensed to ACM.

ACM 978-1-4503-2873-9 /14/09...\$15.00.

<http://dx.doi.org/10.1145/2628136.2628142>

2. Example: Fast Reverse

We shall begin with an example that motivates the rest of the paper: transforming the naïve list reverse function into the so-called “fast reverse” function. This transformation is an instance of the worker/wrapper transformation, and there is an intuitive, informal justification of why this is an optimisation. Here we give this non-rigorous explanation; the remainder of this paper will focus on building the tools to strengthen this to a rigorous argument.

We start with a naïve definition of the reverse function, which takes quadratic time to run as each `append ++` takes time linear in the length of its left argument:

```
reverse      :: [a] → [a]
reverse []   = []
reverse (x : xs) = reverse xs ++ [x]
```

We can write a more efficient version by using a *worker* function *revcat* with a *wrapper* around it that simply applies the worker function with `[]` as the second argument:

```
reverse'     :: [a] → [a]
reverse' xs = revcat xs []
```

The specification for the worker *revcat* is as follows:

```
revcat       :: [a] → [a] → [a]
revcat xs ys = reverse xs ++ ys
```

From this specification we can calculate a new definition that does not depend on *reverse*. Because *reverse* is defined by cases, we will have one calculation for each case.

Case for `[]`:

```
revcat [] ys
= { specification of revcat }
reverse [] ++ ys
= { definition of reverse }
[] ++ ys
= { definition of ++ }
ys
```

Case for `(x : xs)`:

```
revcat (x : xs) ys
= { specification of revcat }
reverse (x : xs) ++ ys
= { definition of reverse }
(reverse xs ++ [x]) ++ ys
= { associativity of ++ }
reverse xs ++ ([x] ++ ys)
= { definition of ++ }
reverse xs ++ (x : ys)
= { specification of revcat }
revcat xs (x : ys)
```

Note the use of associativity of `++` in the third step, which is the only step not simply by definition or specification. Left-associated appends such as `(xs ++ ys) ++ zs` are less time-efficient than the equivalent right-associated appends `xs ++ (ys ++ zs)`, as the former traverses *xs* twice. The intuition here is that the efficiency gain from this step in the proof carries over in some way to the rest of the proof, so that overall our calculated definition of *revcat* is more efficient than its original specification. The calculation gives us the following definition, which runs in linear time:

```
reverse xs      = revcat xs []
revcat [] ys    = ys
revcat (x : xs) ys = revcat xs (x : ys)
```

Unfortunately, there are a number of problems with this approach. Firstly, we calculated *revcat* using the *fold-unfold* style of program calculation [2]. This is an informal calculation, which fails to guarantee total correctness. Thus the resulting *reverse* function may fail in some cases where the original succeeded. Secondly, while we are applying the common pattern of factorising a program into a worker and a wrapper, the reasoning we use is ad-hoc and does not take advantage of this. We would like to abstract out this pattern to make future applications of this technique more straightforward. Finally, while intuitively we can see an efficiency gain from the use of associativity of `++`, this is not a rigorous argument. Put simply, we need rigorous proofs of both *correctness* and *improvement* for our transformation.

3. Worker/Wrapper Transformation

The worker/wrapper transformation, as originally formulated by Gill and Hutton [7], allowed a function written using general recursion to be split into a recursive *worker* function and a *wrapper* function that allows the new definition to be used in the same contexts as the original. The usual application of this technique would be to write the worker to use a different type than the original program that supports more efficient operations, thus hopefully resulting in a more efficient program overall. Gill and Hutton gave conditions for the correctness of the transformation; here we present the more general theory and correctness conditions recently developed by Sculthorpe and Hutton [25].

3.1 The Fix Theory

The idea of the worker/wrapper transformation for fixed-points is as follows. Given a recursive program *prog* of some type *A*, we can write *prog* as some function *f* of itself:

```
prog :: A
prog = f prog
```

We can rewrite this definition so that it is explicitly written using the well-known fixpoint operator *fix*:

```
fix :: (a → a) → a
fix f = f (fix f)
```

resulting in the following definition:

```
prog = fix f
```

Next, we write functions *abs* :: *B* → *A* and *rep* :: *A* → *B* that allow us to convert from the original type *A* to some other type *B* that supports more efficient operations. We finish by constructing a new function *g* : *B* → *B* that allows us to rewrite our original definition of *prog* as follows:

```
prog = abs (fix g)
```

Here *abs* is the *wrapper* function, while *fix g* is the *worker*. The pattern of the worker/wrapper transformation can be captured by a theorem that expresses necessary and sufficient conditions for its correctness [25]. This theorem has assumptions that express the required relationship between the functions *abs* and *rep*, and conditions that provide a specification for the function *g* in terms of *abs*, *rep* and *f*:

Theorem 1 (Worker/Wrapper Factorisation).

Given

$$\begin{array}{ll} \text{abs} : B \rightarrow A & f : A \rightarrow A \\ \text{rep} : A \rightarrow B & g : B \rightarrow B \end{array}$$

satisfying one of the assumptions

$$\begin{array}{ll} (A) \text{ abs} \circ \text{rep} & = \text{id}_A \\ (B) \text{ abs} \circ \text{rep} \circ f & = f \\ (C) \text{ fix } (\text{abs} \circ \text{rep} \circ f) & = \text{fix } f \end{array}$$

and one of the conditions

$$\begin{array}{ll} (1) g = \text{rep} \circ f \circ \text{abs} & (1\beta) \text{ fix } g = \text{fix } (\text{rep} \circ f \circ \text{abs}) \\ (2) g \circ \text{rep} = \text{rep} \circ f & (2\beta) \text{ fix } g = \text{rep } (\text{fix } f) \\ (3) f \circ \text{abs} = \text{abs} \circ g & \end{array}$$

we have the factorisation

$$\text{fix } f = \text{abs } (\text{fix } g)$$

The different assumptions and conditions allow one to choose which will be easiest to verify.

3.2 Proving Fast Reverse Correct

Recall once again the naïve definition of *reverse*:

$$\begin{array}{ll} \text{reverse} & :: [a] \rightarrow [a] \\ \text{reverse } [] & = [] \\ \text{reverse } (x : xs) & = \text{reverse } xs ++ [x] \end{array}$$

As we mentioned before, this naïve implementation is inefficient due to the use of the append operation $++$. We would like to use worker/wrapper factorisation to improve it. The first step is to rewrite the function using *fix*:

$$\begin{array}{ll} \text{reverse} & = \text{fix } \text{rev} \\ \text{rev} & :: ([a] \rightarrow [a]) \rightarrow ([a] \rightarrow [a]) \\ \text{rev } r [] & = [] \\ \text{rev } r (x : xs) & = r \, xs ++ [x] \end{array}$$

The next step in applying worker/wrapper is to select a new type to replace the original type $[a] \rightarrow [a]$, and to write *abs* and *rep* functions to perform the conversions. We can represent a list *xs* by its *difference list* $\lambda ys \rightarrow xs ++ ys$, as first demonstrated by Hughes [12]. Difference lists have the advantage that the usually costly operation of $++$ can be implemented with function composition, typically leading to an increase of efficiency. We write the following functions to convert between the two representations:

$$\begin{array}{ll} \text{type } \text{DiffList } a & = [a] \rightarrow [a] \\ \text{toDiff} & :: [a] \rightarrow \text{DiffList } a \\ \text{toDiff } xs & = \lambda ys \rightarrow xs ++ ys \\ \text{fromDiff} & :: \text{DiffList } a \rightarrow [a] \\ \text{fromDiff } h & = h [] \end{array}$$

We have $\text{fromDiff} \circ \text{toDiff} = \text{id}$:

$$\begin{array}{l} \text{fromDiff } (\text{toDiff } xs) \\ = \{ \text{definition of toDiff} \} \\ \text{fromDiff } (\lambda ys \rightarrow xs ++ ys) \\ = \{ \text{definition of fromDiff} \} \\ (\lambda ys \rightarrow xs ++ ys) [] \\ = \{ \beta\text{-reduction} \} \\ xs ++ [] \\ = \{ [] \text{ is identity of } ++ \} \\ xs \end{array}$$

From these functions it is straightforward to create the actual *abs* and *rep* functions. These convert between the original function type $[a] \rightarrow [a]$ and a new function type $[a] \rightarrow \text{DiffList } a$ where the returned value is represented as a difference list, rather than a regular list:

$$\begin{array}{ll} \text{rep} & :: ([a] \rightarrow [a]) \rightarrow ([a] \rightarrow \text{DiffList } a) \\ \text{rep } h & = \text{toDiff} \circ h \\ \text{abs} & :: ([a] \rightarrow \text{DiffList } a) \rightarrow ([a] \rightarrow [a]) \\ \text{abs } h & = \text{fromDiff} \circ h \end{array}$$

Assumption (A) holds trivially:

$$\begin{array}{l} \text{abs } (\text{rep } h) \\ = \{ \text{definitions of abs and rep} \} \\ \text{fromDiff} \circ \text{toDiff} \circ h \\ = \{ \text{fromDiff} \circ \text{toDiff} = \text{id} \} \\ h \end{array}$$

Now we must verify that the definition of *revcat* that we calculated in the previous section

$$\begin{array}{ll} \text{revcat } [] \, ys & = ys \\ \text{revcat } (x : xs) \, ys & = \text{revcat } xs \, (x : ys) \end{array}$$

satisfies one of the worker/wrapper conditions. We first rewrite *revcat* as an explicit fixed point.

$$\begin{array}{ll} \text{revcat} & = \text{fix } \text{rev}' \\ \text{rev}' \, h \, [] \, ys & = ys \\ \text{rev}' \, h \, (x : xs) \, ys & = h \, xs \, (x : ys) \end{array}$$

We now verify condition (2), $\text{rev}' \circ \text{rep} = \text{rep} \circ \text{rev}$, which expands to $\text{rev}' (\text{rep } r) \, xs = \text{rep } (\text{rev } r) \, xs$. We calculate from the right-hand side, performing case analysis on *xs*. Firstly, we calculate for the case when *xs* is empty:

$$\begin{array}{l} \text{rep } (\text{rev } r) \, [] \\ = \{ \text{definition of rep} \} \\ \text{toDiff } (\text{rev } r \, []) \\ = \{ \text{definition of rev} \} \\ \text{toDiff } [] \\ = \{ \text{definition of toDiff} \} \\ \lambda ys \rightarrow [] ++ ys \\ = \{ [] \text{ is identity of } ++ \} \\ \lambda ys \rightarrow ys \\ = \{ \text{definition of rev}' \} \\ \text{rev}' (\text{rep } r) \, [] \end{array}$$

and then for the case where *xs* is non-empty:

$$\begin{array}{l} \text{rep } (\text{rev } r) \, (x : xs) \\ = \{ \text{definition of rep} \} \\ \text{toDiff } (\text{rev } r \, (x : xs)) \\ = \{ \text{definition of rev} \} \\ \text{toDiff } (r \, xs ++ [x]) \\ = \{ \text{definition of toDiff} \} \\ \lambda ys \rightarrow (r \, xs ++ [x]) ++ ys \\ = \{ \text{associativity and definition of } ++ \} \\ \lambda ys \rightarrow r \, xs ++ (x : ys) \\ = \{ \text{definition of toDiff} \} \\ \lambda ys \rightarrow \text{toDiff } (r \, xs) \, (x : ys) \\ = \{ \text{definition of rep} \} \\ \lambda ys \rightarrow \text{rep } r \, xs \, (x : ys) \\ = \{ \text{definition of rev}' \} \\ \text{rev}' (\text{rep } r) \, (x : xs) \end{array}$$

For total correctness on infinite lists we must also verify the condition holds for the undefined value \perp :

$$\begin{aligned}
& \text{rep } (\text{rev } r) \perp \\
= & \{ \text{definition of } \text{rep} \} \\
& \text{toDiff } (\text{rev } r \perp) \\
= & \{ \text{rev pattern matches on second argument} \} \\
& \text{toDiff } \perp \\
= & \{ \text{definition of } \text{toDiff} \} \\
& \lambda ys \rightarrow \perp \# ys \\
= & \{ \# \text{ strict in first argument} \} \\
& \lambda ys \rightarrow \perp \\
= & \{ \text{rev' pattern matches on second argument} \} \\
& \text{rev' } (\text{rep } r) \perp
\end{aligned}$$

Now that we know our rev' satisfies condition (2), we have a new definition of *reverse*

$$\text{reverse} = \text{abs revcat} = \text{fromDiff} \circ \text{revcat}$$

which eta-expands as follows:

$$\begin{aligned}
\text{reverse } xs &= \text{revcat } xs [] \\
\text{revcat } [] \text{ } ys &= ys \\
\text{revcat } (x : xs) \text{ } ys &= \text{revcat } xs (x : ys)
\end{aligned}$$

The end result is the same improved definition of *reverse* we had before. Thus the worker/wrapper theory has allowed us to formally verify the correctness of our earlier transformation. Furthermore, the use of a general theory has allowed us to avoid the need for induction which would usually be needed to reason about recursive definitions.

4. Improvement Theory

Thus far we have only reasoned about correctness. In order to develop a worker/wrapper theory that can prove *efficiency* properties, we need an operational theory of *program improvement*. More than just expressing extensional information, this should be based on intensional properties of resources that a program requires. For the purpose of this paper, the resource we shall consider is execution time.

We have two main design goals for our operational theory. Firstly, it ought to be based on the operational semantics of a realistic programming language, so that conclusions we draw from it are as applicable as possible. Secondly, it should be amenable to techniques such as (in)equational reasoning, as these are the techniques we used to apply the worker/wrapper correctness theory.

For the first goal, we use a language with similar syntax and semantics to GHC Core, except that arguments to functions are required to be atomic, as was the case in earlier versions of the language [20]. (Normalisation of the current version of GHC Core into this form is straightforward.) The language is call-by-need, reflecting the use of lazy evaluation in Haskell. The efficiency behaviour of call-by-need programs is notoriously counterintuitive. Our hope is that providing formal techniques for reasoning about call-by-need efficiency we will go some way toward easing this problem.

For the second goal, our theory must be based around relation R that is a preorder, as transitivity and reflexivity are necessary for inequational reasoning to be valid. Furthermore, to support reasoning in a compositional manner, it is essential to allow substitution. That is, given terms M and N , if $M R N$ then $\mathbb{C}[M] R \mathbb{C}[N]$ should also hold for any context \mathbb{C} . A relation R that satisfies both of these properties is called a *precongruence*.

A naïve approach to measuring execution time would be to simply count the number of steps taken to evaluate a term to some normal form, and consider that a term M

is more efficient than a term N if its evaluation finishes in fewer steps. The resulting relation is clearly a preorder; however it is not a precongruence in a call-by-need setting, because meaningful computations can be done with terms that are not fully normalised. For example, just because M normalises and N does not, it does not follow that M is necessarily more efficient in *all* contexts.

The approach we use is due to Moran and Sands [15]. Rather than counting the steps taken to normalise a term, we compare the steps taken in *all* contexts, and only say that M is improved by N if for any context \mathbb{C} , the term $\mathbb{C}[M]$ requires no more evaluation steps than the term $\mathbb{C}[N]$. The result is a relation that is trivially a precongruence: it inherits transitivity and reflexivity from the numerical ordering \leq , and is substitutive by definition.

Improvement theory [23] was originally developed for call-by-name languages by Sands [21]. The remainder of this section presents the call-by-need time improvement theory due to Moran and Sands [15], which will provide the setting for our operational worker/wrapper theory. The essential difference between call-by-name and call-by-need is that the latter implements a *sharing* strategy, avoiding the repeated evaluation of terms that are used more than once.

4.1 Operational Semantics of the Core Language

We shall begin by presenting the operational model that forms the basis of this improvement theory. The semantics presented here are originally due to Sestoft [27].

We start from a set of variables Var and a set of constructors Con . We assume all constructors have a fixed arity. The grammar of terms is as follows:

$$\begin{aligned}
x, y, z &\in Var \\
c &\in Con \\
M, N &::= x \\
&| \lambda x \rightarrow M \\
&| M \ x \\
&| \text{let } \{\vec{x} = \vec{M}\} \text{ in } N \\
&| c \ \vec{x} \\
&| \text{case } M \text{ of } \{c_i \ \vec{x}_i \rightarrow N_i\}
\end{aligned}$$

We use $\vec{x} = \vec{M}$ as a shorthand for a list of bindings of the form $x = M$. Similarly, we use $c_i \ \vec{x}_i \rightarrow N_i$ as a shorthand for a list of cases of the form $c \ \vec{x} \rightarrow N$. All constructors are assumed to be saturated, that is, we assume that any \vec{x} that is the operand of a constructor c has length equal to the arity of c . Literals are represented by constructors of arity 0. We treat α -equivalent terms as identical.

A term is a *value* if it is of the form $c \ \vec{x}$ or $\lambda x \rightarrow M$. In Haskell this is referred to as a *weak head normal form*. We shall use letters such as V, W to denote value terms.

Term contexts take the following form, with substitution defined in the obvious way.

$$\begin{aligned}
\mathbb{C}, \mathbb{D} &::= [-] \\
&| x \\
&| \lambda x \rightarrow \mathbb{C} \\
&| \mathbb{C} \ x \\
&| \text{let } \{\vec{x} = \vec{\mathbb{C}}\} \text{ in } \mathbb{D} \\
&| c \ \vec{x} \\
&| \text{case } \mathbb{C} \text{ of } \{c_i \ \vec{x}_i \rightarrow \mathbb{D}_i\}
\end{aligned}$$

A value context is a context that is either a lambda abstraction or a constructor applied to variables.

The restriction that the arguments of functions and constructors always be variables has the effect that all bindings

$\langle \Gamma \{x = M\}, x, S \rangle$	$\rightarrow \langle \Gamma, M, \#x : S \rangle$	{ LOOKUP }
$\langle \Gamma, V, \#x : S \rangle$	$\rightarrow \langle \Gamma \{x = V\}, V, S \rangle$	{ UPDATE }
$\langle \Gamma, M x, S \rangle$	$\rightarrow \langle \Gamma, M, x : S \rangle$	{ UNWIND }
$\langle \Gamma, \lambda x \rightarrow M, y : S \rangle$	$\rightarrow \langle \Gamma, M [y / x], S \rangle$	{ SUBST }
$\langle \Gamma, \text{case } M \text{ of } \text{alts}, S \rangle$	$\rightarrow \langle \Gamma, M, \text{alts} : S \rangle$	{ CASE }
$\langle \Gamma, c_j \vec{y}, \{c_i \vec{x}_i \rightarrow N_i\} : S \rangle$	$\rightarrow \langle \Gamma, N_j [\vec{y} / \vec{x}_j], S \rangle$	{ BRANCH }
$\langle \Gamma, \text{let } \{\vec{x} = \vec{M}\} \text{ in } N, S \rangle$	$\rightarrow \langle \Gamma \{\vec{x} = \vec{M}\}, N, S \rangle$	{ LETREC }

Figure 1. The call-by-need abstract machine

made during evaluation must have been created by a **let**. Sometimes we will use $M N$ (where N is not a variable) as a shorthand for **let** $\{x = N\}$ **in** $M x$, where x is fresh. We use this shorthand for both terms and contexts.

An abstract machine for executing terms in the language maintains a state $\langle \Gamma, M, S \rangle$ consisting of: a heap Γ , given by a set of bindings from variables to terms; the term M currently being evaluated; the evaluation stack S , given by a list of tokens used by the abstract machine. The machine works by evaluating the current term to a value, and then decides what to do with the value based on the top of the stack. Bindings generated by **let** constructs are put on the heap, and only taken off when performing a LOOKUP. A LOOKUP executes by putting a token on the stack representing where the term was looked up, and then evaluating that term to value form before replacing it on the heap. In this way, each binding is only ever evaluated at most once. The semantics of the machine is given in Figure 1. Note that the LETREC rule assumes that \vec{x} is disjoint from the domain of Γ ; if not, we need only α -rename so that this is the case.

4.2 The Cost Model and Improvement Relations

Now that we have a semantics for our model, we must devise a *cost model* for this semantics. The natural way to do this for an operational semantics is to count steps taken to evaluate a given term. We use the notation $M \downarrow^n$ to mean the abstract machine progresses from the initial state $\langle \emptyset, M, \epsilon \rangle$ to some final state $\langle \Gamma, V, \epsilon \rangle$ with n occurrences of the LOOKUP step. It is sufficient to count LOOKUP steps because the total number of steps is bounded by a linear function of the number of LOOKUP steps [15]. Furthermore, we use the notation $M \downarrow^{\leq n}$ to mean that $M \downarrow^m$ for some $m \leq n$.

From this, we can define our improvement relation. We say that “ M is improved by N ”, written $M \preceq N$, if the following statement holds for all contexts \mathbb{C} :

$$\mathbb{C}[M] \downarrow^m \implies \mathbb{C}[N] \downarrow^{\leq m}$$

In other words, a term M is improved by a term N if N takes no more steps to evaluate than M in all contexts. That this relation is a congruence follows immediately from the definition, and that it is a preorder follows from the fact that \leq is itself a preorder. We sometimes write $M \preceq N$ for $N \preceq M$. If both $M \preceq N$ and $M \preceq N$, we write $M \approx N$ and say that M and N are *cost-equivalent*.

For convenience, we define a “tick” operation on terms that adds exactly one unit of cost to a term:

$$\check{M} \equiv \text{let } \{x = M\} \text{ in } x \quad \{ \text{where } x \text{ is free in } M \}$$

This definition for \check{M} takes exactly two steps to evaluate to M : one to add the binding to the heap, and the other to look it up. Only one of these steps is a LOOKUP step, so the result is that the cost of evaluating the term is increased by exactly one. Using ticks allows us to annotate terms with in-

dividual units of cost, allowing us to use rules to “push” cost around a term, making the calculations more convenient. We could also define the tick operation by adding it to the grammar of terms and modifying the abstract machine and cost model accordingly, but this definition is equivalent. We have the following law: $\check{M} \preceq M$.

The improvement relation \preceq covers when one term is at least as efficient as another in all contexts, but this is a very strong statement. We use the notion of “weak improvement” when one term is at least as efficient as another within a constant factor. Specifically, we say M is weakly improved by N , written $M \preceq_w N$, if there exists a linear function $f(x) = kx + c$ (where $k, c \geq 0$) such that the following statement holds for all contexts \mathbb{C} :

$$\mathbb{C}[M] \downarrow^m \implies \mathbb{C}[N] \downarrow^{\leq f(m)}$$

This can be read as “replacing M with N may make programs worse, but cannot make them *asymptotically* worse”. We use symbols \preceq and \approx for inverse and equivalence analogously as for standard improvement.

Because weak improvement ignores constant factors, we have the following *tick introduction/elimination* law:

$$M \approx \check{M}$$

It follows from this that any improvement $M \preceq N$ can be weakened to a weak improvement $M' \preceq_w N'$ where M' and N' denote the terms M and N with all the ticks removed.

The last notation we define is *entailment*, which is used when we have a chain of improvements that all apply with respect to a particular set of definitions. Specifically, where $\Gamma = \{\vec{x} = \vec{V}\}$ is a list of bindings, we write:

$$\Gamma \vdash M_1 \preceq M_2 \preceq \dots \preceq M_n$$

to mean:

$$\text{let } \Gamma \text{ in } M_1 \preceq \text{let } \Gamma \text{ in } M_2 \preceq \dots \preceq \text{let } \Gamma \text{ in } M_n$$

4.3 Selected Laws

We finish this section with a selection of laws taken from [15]. The first two are β -reduction rules. The following cost equivalence holds for function application:

$$(\lambda x \rightarrow M) y \approx M [y / x]$$

This holds because the abstract machine evaluates the left-hand-side to the right-hand-side without performing any LOOKUPS, resulting the same heap and stack as before. Note that the substitution is variable-for-variable, as the grammar for our language requires that the argument to function application always be a variable.

In general, where a term M can be evaluated to a term M' , we have the following relationships:

$$\begin{aligned} M &\preceq M' \\ M' &\approx M \end{aligned}$$

The latter fact may be non-obvious, but it holds because evaluating a term will produce a constant number of ticks, and tick-elimination is a weak cost-equivalence. In this manner we can see that partial evaluation by itself will never save more than a constant-factor of time.

The following cost equivalence allows us to substitute a variable for its binding. However, note that this is only valid for *values*, as bindings to other terms will be modified in the course of execution. We thus call this rule *value-β*.

$$\begin{aligned} & \text{let } \{x = V, \vec{y} = \vec{C}[x]\} \text{ in } \mathbb{D}[x] \\ & \Downarrow \\ & \text{let } \{x = V, \vec{y} = \vec{C}[\sqrt{V}]\} \text{ in } \mathbb{D}[\sqrt{V}] \end{aligned}$$

The following law allows us to move let bindings in and out of a context when the binding is to a value. Note that we assume that x does not appear free in \mathbb{C} , which can be ensured by α -renaming, and that no free variables in V are captured in \mathbb{C} . We call this rule *value let-floating*.

$$\mathbb{C}[\text{let } \{x = V\} \text{ in } M] \Downarrow \text{let } \{x = V\} \text{ in } \mathbb{C}[M]$$

We also have a *garbage collection* law allowing us to remove unused bindings. Assuming that x is not free in \vec{N} or L , we have the following cost equivalence:

$$\text{let } \{x = M; \vec{y} = \vec{N}\} \text{ in } L \Downarrow \text{let } \{\vec{y} = \vec{N}\} \text{ in } L$$

The final law we present here is the rule of *improvement induction*. The version that we present is stronger than the version in [15], but can be obtained by a simple modification of the proof given there. For any set of value bindings Γ and context \mathbb{C} , we have the following rule:

$$\frac{\Gamma \vdash M \Downarrow \sqrt{\mathbb{C}[M]} \quad \Gamma \vdash \sqrt{\mathbb{C}[N]} \Downarrow N}{\Gamma \vdash M \Downarrow N}$$

This allows us to prove an $M \Downarrow N$ simply by finding a context \mathbb{C} where we can “unfold” M to $\sqrt{\mathbb{C}[M]}$ and “fold” $\sqrt{\mathbb{C}[N]}$ to N . In other words, the following proof is valid:

$$\begin{aligned} & \Gamma \vdash M \\ & \Downarrow \\ & \sqrt{\mathbb{C}[M]} \\ & \Downarrow \{ \text{hypothesis} \} \\ & \sqrt{\mathbb{C}[N]} \\ & \Downarrow \\ & N \end{aligned}$$

In this way the technique is similar to proof principles such as guarded coinduction [4, 28].

As a corollary to this law, we have the following law for *cost-equivalence* improvement induction. For any set of value bindings Γ and context \mathbb{C} , we have:

$$\frac{\Gamma \vdash M \Downarrow \sqrt{\mathbb{C}[M]} \quad \Gamma \vdash \sqrt{\mathbb{C}[N]} \Downarrow N}{\Gamma \vdash M \Downarrow N}$$

The proof is simply to start from the assumptions and make two applications of improvement induction: first to prove $M \Downarrow N$, and second to prove $N \Downarrow M$.

5. Worker/Wrapper and Improvement

In this section, we prove a factorisation theorem for improvement theory analogous to the worker/wrapper factorisation theorem given in section 3.1. Before we do this, however, we must prove two preliminary results: a *rolling rule* and a *fusion rule*. Rolling and fusion are central to the worker/wrapper transformation [7, 13], so it is only natural that we would need versions of these to apply worker/wrapper transformation in this context.

5.1 Preliminary Results

The first rule we prove is the *rolling rule*, so named because of its similarity to the rolling rule for least-fixed points. In particular, for any pair of value contexts \mathbb{F}, \mathbb{G} , we have the following weak cost equivalence:

$$\text{let } \{x = \mathbb{F}[\mathbb{G}[x]]\} \text{ in } \mathbb{G}[x] \Downarrow \text{let } \{x = \mathbb{G}[\mathbb{F}[x]]\} \text{ in } x$$

The proof begins with an application of cost-equivalence improvement induction. We let $\Gamma = \{x = \mathbb{F}[\sqrt{\mathbb{G}[x]}], y = \mathbb{G}[\sqrt{\mathbb{F}[y]}], M = \sqrt{\mathbb{G}[x]}, N = y, \mathbb{C} = \mathbb{G}[\sqrt{\mathbb{F}[-]}]$. The premises of induction are proved as follows:

$$\begin{aligned} \Gamma \vdash M & \equiv \{ \text{definitions} \} \\ & \sqrt{\mathbb{G}[x]} \\ & \Downarrow \{ \text{value-}\beta \} \\ & \sqrt{\mathbb{G}[\sqrt{\mathbb{F}[\sqrt{\mathbb{G}[x]}]]} \\ & \equiv \{ \text{definitions} \} \\ & \sqrt{\mathbb{C}[M]} \end{aligned}$$

and

$$\begin{aligned} \Gamma \vdash \sqrt{\mathbb{C}[N]} & \equiv \{ \text{definitions} \} \\ & \sqrt{\mathbb{G}[\sqrt{\mathbb{F}[y]}]} \\ & \Downarrow \{ \text{value-}\beta \} \\ & y \\ & \equiv \{ \text{definitions} \} \\ & N \end{aligned}$$

Thus we can conclude $\Gamma \vdash M \Downarrow N$, or equivalently $\text{let } \Gamma \text{ in } M \Downarrow \text{let } \Gamma \text{ in } N$. We expand this out and apply garbage collection to remove the unused bindings:

$$\text{let } \{x = \mathbb{F}[\sqrt{\mathbb{G}[x]}]\} \text{ in } \sqrt{\mathbb{G}[x]} \Downarrow \text{let } \{y = \mathbb{G}[\sqrt{\mathbb{F}[y]}]\} \text{ in } y$$

By applying α -renaming and weakening we obtain the desired result. The second rule we prove is *letrec-fusion*, analogous to fixed-point fusion. For any value contexts \mathbb{F}, \mathbb{G} , we have the following implication:

$$\begin{aligned} & \mathbb{H}[\sqrt{\mathbb{F}[x]}] \Downarrow \mathbb{G}[\sqrt{\mathbb{H}[x]}] \\ & \Rightarrow \\ & \text{let } \{x = \mathbb{F}[x]\} \text{ in } \mathbb{H}[x] \Downarrow \text{let } \{x = \mathbb{G}[x]\} \text{ in } x \end{aligned}$$

For the proof, we assume the premise and proceed by improvement induction. Let $\Gamma = \{x = \mathbb{F}[x], y = \mathbb{G}[y]\}$, $M = \sqrt{\mathbb{H}[x]}$, $N = y$, $\mathbb{C} = \mathbb{G}$. The premises are proved by:

$$\begin{aligned} \Gamma \vdash M & \equiv \{ \text{by definitions} \} \\ & \sqrt{\mathbb{H}[x]} \\ & \Downarrow \{ \text{value beta} \} \\ & \sqrt{\mathbb{H}[\sqrt{\mathbb{F}[x]}]} \\ & \Downarrow \{ \text{by assumption} \} \\ & \sqrt{\mathbb{G}[\sqrt{\mathbb{H}[x]}]} \\ & \equiv \{ \text{definition} \} \\ & \sqrt{\mathbb{C}[M]} \end{aligned}$$

and

$$\begin{aligned} \Gamma \vdash \sqrt{\mathbb{C}[N]} & \equiv \{ \text{by definitions} \} \\ & \sqrt{\mathbb{G}[y]} \\ & \Downarrow \{ \text{value beta} \} \\ & y \\ & \equiv \{ \text{definition} \} \\ & N \end{aligned}$$

Thus we conclude that $\Gamma \vdash M \preceq N$. Expanding and applying garbage collection, we obtain the following:

$\text{let } \{x = \mathbb{F}[x]\} \text{ in } \mathbb{H}[x] \preceq \text{let } y = \mathbb{G}[y] \text{ in } y$

Again we obtain the desired result via weakening and α -renaming. As improvement induction is symmetrical, we can also prove the following dual fusion law, in which the improvement relations are reversed:

$\mathbb{H}[\sqrt{\mathbb{F}[x]}] \preceq \mathbb{G}[\sqrt{\mathbb{H}[x]}]$
 \Rightarrow
 $\text{let } \{x = \mathbb{F}[x]\} \text{ in } \mathbb{H}[x] \preceq \text{let } \{x = \mathbb{G}[x]\} \text{ in } x$

For both the rolling and fusion rules, we first proved a version of the conclusion with normal improvement, and then weakened to weak improvement. We do this to avoid having to deal with ticks, and because the weaker version is strong enough for our purposes.

Moran and Sands also prove their own fusion law. This law requires that the context \mathbb{H} satisfy a form of *strictness*. Specifically, For any value contexts \mathbb{F} , \mathbb{G} and fresh variable x , we have the following implication:

$\mathbb{H}[\mathbb{F}[x]] \preceq \mathbb{G}[\mathbb{H}[x]] \wedge \text{strict } (\mathbb{H})$
 \Rightarrow
 $\text{let } \{x = \mathbb{F}[x]\} \text{ in } \mathbb{C}[\mathbb{H}[x]] \preceq \text{let } \{x = \mathbb{G}[x]\} \text{ in } \mathbb{C}[x]$

This version of fusion has the advantage of having a stronger conclusion, but its strictness side-condition and lack of symmetry make it unsuitable for our purposes.

5.2 The Worker/Wrapper Improvement Theorem

Using the above set of rules, we can prove the following *worker/wrapper improvement* theorem, giving conditions under which a program factorisation is a time improvement:

Theorem 2 (Worker/Wrapper Improvement).

Given value contexts Abs , Rep , \mathbb{F} , \mathbb{G} for which x is free satisfying one of the assumptions

- (A) $\text{Abs}[\text{Rep}[x]] \preceq x$
- (B) $\text{Abs}[\text{Rep}[\mathbb{F}[x]]] \preceq \mathbb{F}[x]$
- (C) $\text{let } x = \text{Abs}[\text{Rep}[\mathbb{F}[x]]] \text{ in } x \preceq \text{let } x = \mathbb{F}[x] \text{ in } x$

and one of the conditions

- (1) $\mathbb{G}[x] \preceq \text{Rep}[\mathbb{F}[\text{Abs}[x]]]$
- (2) $\mathbb{G}[\sqrt{\text{Rep}[x]}] \preceq \text{Rep}[\sqrt{\mathbb{F}[x]}]$
- (3) $\text{Abs}[\sqrt{\mathbb{G}[x]}] \preceq \mathbb{F}[\sqrt{\text{Abs}[x]}]$
- (1 β) $\text{let } x = \mathbb{G}[x] \text{ in } x \preceq \text{let } x = \text{Rep}[\mathbb{F}[\text{Abs}[x]]] \text{ in } x$
- (2 β) $\text{let } x = \mathbb{G}[x] \text{ in } x \preceq \text{let } x = \mathbb{F}[x] \text{ in } \text{Rep}[x]$

we have the improvement

$\text{let } x = \mathbb{F}[x] \text{ in } x \preceq \text{let } x = \mathbb{G}[x] \text{ in } \text{Abs}[x]$

Given a recursive program $\text{let } x = \mathbb{F}[x] \text{ in } x$ and *abstraction* and *representation* contexts Abs and Rep , this theorem gives us conditions we can use to derive a factorised program $\text{let } x = \mathbb{G}[x] \text{ in } \text{Abs}[x]$. This factorised program will be at worst a constant factor slower than the original program, but can potentially be asymptotically faster. In other words, we have conditions that guarantee that such an optimisation is “safe” with respect to time performance.

The proof given in [25] for the original factorisation theorem centers on the use of the rolling and fusion rules. Because we have proven analogous rules in our setting, the proofs can be adapted fairly straightforwardly, simply by keeping the general form of the proofs and using the rules

of improvement theory as structural rules that fit between the original steps. The details are as follows.

We begin by noting that $(A) \Rightarrow (B) \Rightarrow (C)$, as in the original case. The first implication $(A) \Rightarrow (B)$ no longer follows immediately, but the proof is simple. Letting y be a fresh variable, we reason as follows:

$\text{Abs}[\text{Rep}[\mathbb{F}[y]]]$
 $\preceq \{ \text{garbage collection, value-}\beta \}$
 $\preceq \text{let } x = \mathbb{F}[y] \text{ in } \text{Abs}[\text{Rep}[x]]$
 $\preceq \{ (A) \}$
 $\preceq \text{let } x = \mathbb{F}[y] \text{ in } x$
 $\preceq \{ \text{value-}\beta, \text{garbage collection} \}$
 $\preceq \mathbb{F}[y]$

The final step is to observe that as both x and y are fresh, we can substitute one for the other and the relationship between the terms will remain the same. Hence, we can conclude (B) .

As in the original theorem, we have that (1) implies (1 β) by simple application of substitution, (2) implies (2 β) by fusion and (3) implies the conclusion also by fusion. Under assumption (C) , we have that (1 β) and (2 β) are equivalent. We show this by proving their right hand sides cost-equivalent, after which we can simply apply transitivity.

$\text{let } x = \mathbb{F}[x] \text{ in } \text{Rep}[x]$
 $\preceq \{ \text{value-}\beta \}$
 $\preceq \text{let } x = \mathbb{F}[x] \text{ in } \text{Rep}[\mathbb{F}[x]]$
 $\preceq \{ \text{value let-floating} \}$
 $\preceq \text{Rep}[\mathbb{F}[\text{let } x = \mathbb{F}[x] \text{ in } x]]$
 $\preceq \{ (C) \}$
 $\preceq \text{Rep}[\mathbb{F}[\text{let } x = \text{Abs}[\text{Rep}[\mathbb{F}[x]]] \text{ in } x]]$
 $\preceq \{ \text{value let-floating} \}$
 $\preceq \text{let } x = \text{Abs}[\text{Rep}[\mathbb{F}[x]]] \text{ in } \text{Rep}[\mathbb{F}[x]]$
 $\preceq \{ \text{rolling} \}$
 $\preceq \text{let } x = \text{Rep}[\mathbb{F}[\text{Abs}[x]]] \text{ in } x$

Finally, we must show that condition (1 β) and assumption (C) together imply the conclusion. This follows exactly the same pattern of reasoning as the original proof, with the addition of two applications of value-let floating:

$\text{let } x = \mathbb{F}[x] \text{ in } x$
 $\preceq \{ (C) \}$
 $\preceq \text{let } x = \text{Abs}[\text{Rep}[\mathbb{F}[x]]] \text{ in } x$
 $\preceq \{ \text{rolling} \}$
 $\preceq \text{let } x = \text{Rep}[\mathbb{F}[\text{Abs}[x]]] \text{ in } \text{Abs}[x]$
 $\preceq \{ \text{value let-floating} \}$
 $\preceq \text{Abs}[\text{let } x = \text{Rep}[\mathbb{F}[\text{Abs}[x]]] \text{ in } x]$
 $\preceq \{ (1\beta) \}$
 $\preceq \text{Abs}[\text{let } x = \mathbb{G}[x] \text{ in } x]$
 $\preceq \{ \text{value let-floating} \}$
 $\preceq \text{let } x = \mathbb{G}[x] \text{ in } \text{Abs}[x]$

We conclude this section by discussing a few important points about the worker/wrapper improvement theorem and its applications. Firstly, we note that the condition (A) will never actually hold. To see this, we let Ω be a divergent term; that is, one that the abstract machine will never finish evaluating. By substituting into the context $\text{let } x = \Omega \text{ in } [-]$, we obtain the following cost-equivalence:

$\text{let } x = \Omega \text{ in } \text{Abs}[\text{Rep}[x]] \not\preceq \text{let } x = \Omega \text{ in } x$

This is clearly false, as the left-hand side will terminate almost immediately (as Abs is a value context), while the right-hand side will diverge. Thus we see that assumption (A) is impossible to satisfy. We leave it in the theorem

for completeness of the analogy with the earlier theorem from section 3.1. In situations where (A) would have been used with the earlier theory, the weaker assumption (B) can always be used instead. As we will see later with the examples, frequently only very few properties of the context \mathbb{F} will be used in the proof of (B). A *typed* improvement theory might allow these properties to be assumed of x instead, thus making (A) useful again.

Secondly, we note the restriction to value contexts. This is not actually a particularly severe restriction: for the common application of recursively-defined functions, it is fairly straightforward to ensure that all contexts be of the form $\lambda x \rightarrow \mathbb{C}$. For other applications it may be more difficult to find **Abs** and **Rep** contexts with the required relationship.

Finally, we note that only conditions (2) and (3) use normal improvement, with all other assumptions and conditions using the weaker version. This is because weak improvement is not strong enough to permit the use of fusion, which these conditions rely on. This makes these conditions harder to prove. However, when these conditions are used, their strength allows us to narrow down the source of any constant-factor slowdown that may take place.

6. Examples

6.1 Reversing a List

In this section we shall demonstrate the utility of our theory with two practical examples. We begin by revisiting the earlier example of reversing a list. In order to apply our theory, we must first write *reverse* as a recursive let:

$$\begin{aligned} \text{reverse} &= \text{let } \{f = \text{Revbody } [f]\} \text{ in } f \\ \text{Revbody}[M] &= \lambda xs \rightarrow \text{case } xs \text{ of} \\ &\quad [] \rightarrow [] \\ &\quad (y : ys) \rightarrow M \text{ } ys ++ [y] \end{aligned}$$

The *abs* and *rep* functions from before give rise to the following contexts:

$$\begin{aligned} \text{Abs}[M] &= \lambda xs \rightarrow M \text{ } xs \ [] \\ \text{Rep}[M] &= \lambda xs \rightarrow \lambda ys \rightarrow M \text{ } xs ++ ys \end{aligned}$$

We also require some extra theoretical machinery that we have yet to introduce. To start with, we must assume some rules about the append operation $++$. The following associativity rules were proved by Moran and Sands [15].

$$\begin{aligned} (xs ++ ys) ++ zs &\preceq xs ++ (ys ++ zs) \\ xs ++ (ys ++ zs) &\preceq (xs ++ ys) ++ zs \end{aligned}$$

We assume the following identity improvement as well, which follows from theorems also proved in [15]:

$$[] ++ xs \preceq xs$$

We also require the notion of an *evaluation context*. An evaluation context is a context where evaluation is impossible unless the hole is filled, and have the following form:

$$\begin{aligned} \mathbb{E} &:: \mathbb{A} \\ &| \text{let } \{\vec{x} = \vec{M}\} \text{ in } \mathbb{A} \\ &| \text{let } \{\vec{y} = \vec{M}; \\ &\quad x_0 = \mathbb{A}_0[x_1]; \\ &\quad x_1 = \mathbb{A}_1[x_2]; \\ &\quad \dots \\ &\quad x_n = \mathbb{A}_n\} \\ &\text{in } \mathbb{A}[x_0] \end{aligned}$$

$$\begin{aligned} \mathbb{A} &:: [-] \\ &| \mathbb{A} \ x \\ &| \text{case } \mathbb{A} \text{ of } \{c_i \ \vec{x}_i \rightarrow M_i\} \end{aligned}$$

Note that a context of this form must have exactly one hole. The usefulness of evaluation contexts is that they satisfy some special laws. We use the following in this example:

$$\begin{aligned} &\mathbb{E}[\sqrt{M}] \\ &\Downarrow \{ \text{tick floating} \} \\ &\sim \sqrt{\mathbb{E}[M]} \\ &\mathbb{E}[\text{case } M \text{ of } \{c_i \ \vec{x}_i \rightarrow N_i\}] \\ &\Downarrow \{ \text{case floating} \} \\ &\sim \text{case } M \text{ of } \{c_i \ \vec{x}_i \rightarrow \mathbb{E}[N_i]\} \\ &\mathbb{E}[\text{let } \{\vec{x} = \vec{M}\} \text{ in } N] \\ &\Downarrow \{ \text{let floating} \} \\ &\sim \text{let } \{\vec{x} = \vec{M}\} \text{ in } \mathbb{E}[N] \end{aligned}$$

We conclude by noting that while the context $[-] ++ ys$ is not strictly speaking an evaluation context (as the hole is in the wrong place), it is cost-equivalent to an evaluation context and so also satisfies these laws. The proof is as follows:

$$\begin{aligned} &[-] ++ ys \\ &\equiv \{ \text{desugaring} \} \\ &\quad (\text{let } \{xs = [-]\} \text{ in } (++) \ xs) \ ys \\ &\Downarrow \{ \text{let floating } [-] \ ys \} \\ &\sim \text{let } \{xs = [-]\} \text{ in } (++) \ xs \ ys \\ &\Downarrow \{ \text{unfolding } ++ \} \\ &\sim \text{let } \{xs = [-]\} \text{ in} \\ &\quad \sqrt{\text{case } xs \text{ of} \\ &\quad \quad [] \rightarrow ys \\ &\quad \quad (z : zs) \rightarrow \text{let } \{rs = (++) \ zs \ ys\} \text{ in } z : rs} \\ &\Downarrow \{ \text{desugaring tick and collecting lets} \} \\ &\sim \text{let } \{xs = [-]; \\ &\quad r = \text{case } xs \text{ of} \\ &\quad \quad [] \rightarrow ys \\ &\quad \quad (z : zs) \rightarrow \text{let } \{rs = (++) \ zs \ ys\} \text{ in } z : rs\} \\ &\quad \text{in } r \end{aligned}$$

Now we can begin the example proper. We start by verifying that **Abs** and **Rep** satisfy one of the worker/wrapper assumptions. While earlier we used (A) for this example, the corresponding assumption for worker/wrapper improvement is unsatisfiable. Thus we instead verify assumption (B). The proof is fairly straightforward:

$$\begin{aligned} &\text{Abs}[\text{Rep}[\text{Revbody}[f]]] \\ &\equiv \{ \text{definitions} \} \\ &\quad \lambda xs \rightarrow (\lambda xs \rightarrow \lambda ys \rightarrow \text{Revbody}[f] \ xs ++ ys) \ xs \ [] \\ &\Downarrow \{ \beta\text{-reduction} \} \\ &\sim \lambda xs \rightarrow \text{Revbody}[f] \ xs ++ [] \\ &\equiv \{ \text{definition of Revbody} \} \\ &\quad \lambda xs \rightarrow (\lambda xs \rightarrow \text{case } xs \text{ of} \\ &\quad \quad [] \rightarrow [] \\ &\quad \quad (y : ys) \rightarrow f \ ys ++ [y]) \ xs ++ [] \\ &\Downarrow \{ \beta\text{-reduction} \} \\ &\sim \lambda xs \rightarrow (\text{case } xs \text{ of} \\ &\quad \quad [] \rightarrow [] \\ &\quad \quad (y : ys) \rightarrow f \ ys ++ [y]) ++ [] \\ &\Downarrow \{ \text{case floating } [-] ++ [] \} \\ &\sim \lambda xs \rightarrow \text{case } xs \text{ of} \\ &\quad \quad [] \rightarrow [] ++ [] \\ &\quad \quad (y : ys) \rightarrow (f \ ys ++ [y]) ++ [] \end{aligned}$$

$$\begin{aligned}
& \Downarrow \{ \text{associativity is weak cost equivalence} \} \\
& \approx \lambda xs \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad [] \rightarrow [] \mathbin{++} [] \\
& \quad (y : ys) \rightarrow f \, ys \mathbin{++} ([y] \mathbin{++} []) \\
& \Downarrow \{ \text{evaluating } [] \mathbin{++} [], [y] \mathbin{++} [] \} \\
& \approx \lambda xs \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad [] \rightarrow [] \\
& \quad (y : ys) \rightarrow f \, ys \mathbin{++} [y] \\
& \equiv \{ \text{definition of revbody} \} \\
& \quad \mathbf{Revbody} \, [f]
\end{aligned}$$

As before, we use condition (2) to derive our \mathbb{G} . The derivation is somewhat more involved than before, requiring some care with the manipulation of ticks.

$$\begin{aligned}
& \mathbf{Rep}[\checkmark \mathbf{Revbody}[f]] \\
& \equiv \{ \text{definitions} \} \\
& \quad \lambda xs \rightarrow \lambda ys \rightarrow \\
& \quad \quad (\checkmark \lambda xs \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad \quad [] \rightarrow [] \\
& \quad \quad \quad (z : zs) \rightarrow f \, zs \mathbin{++} [z]) \, xs \mathbin{++} ys \\
& \Downarrow \{ \text{float tick out of } [-] \, xs \mathbin{++} ys \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \\
& \quad \checkmark((\lambda xs \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow [] \\
& \quad \quad (z : zs) \rightarrow f \, zs \mathbin{++} [z]) \, xs \mathbin{++} ys) \\
& \Downarrow \{ \beta\text{-reduction} \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \checkmark((\mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow [] \\
& \quad \quad (z : zs) \rightarrow f \, zs \mathbin{++} [z]) \mathbin{++} ys) \\
& \Downarrow \{ \text{case floating } [-] \mathbin{++} ys \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \checkmark(\mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow [] \mathbin{++} ys \\
& \quad \quad (z : zs) \rightarrow (f \, zs \mathbin{++} [z]) \mathbin{++} ys) \\
& \Downarrow \{ \text{associativity and identity of } ++ \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \checkmark(\mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow f \, zs \mathbin{++} ([z] \mathbin{++} ys)) \\
& \Downarrow \{ \text{evaluating } [y] \mathbin{++} ys \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \checkmark(\mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow f \, zs \mathbin{++} (z : ys)) \\
& \Downarrow \{ \text{case floating tick } (\star) \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow \checkmark \, ys \\
& \quad \quad (z : zs) \rightarrow \checkmark(f \, zs \mathbin{++} (z : ys)) \\
& \Downarrow \{ \text{removing a tick} \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow \checkmark(f \, zs \mathbin{++} (z : ys)) \\
& \Downarrow \{ \text{desugaring} \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow \\
& \quad \quad \quad \checkmark(\mathbf{let} \, ws = (z : ys) \, \mathbf{in} \\
& \quad \quad \quad \quad f \, zs \mathbin{++} ws) \\
& \Downarrow \{ \beta\text{-expansion} \} \\
& \approx \lambda xs \rightarrow \lambda ys \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow \\
& \quad \quad \quad \checkmark \mathbf{let} \, ws = (z : ys) \, \mathbf{in} \\
& \quad \quad \quad \quad (\lambda as \rightarrow \lambda bs \rightarrow f \, as \mathbin{++} bs) \, zs \, ws \\
& \Downarrow \{ \text{tick floating } [-] \, zs \, ws \}
\end{aligned}$$

$$\begin{aligned}
& \lambda xs \rightarrow \lambda ys \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow \\
& \quad \quad \quad \mathbf{let} \, ws = (z : ys) \, \mathbf{in} \\
& \quad \quad \quad \quad (\checkmark \lambda as \rightarrow \lambda bs \rightarrow f \, as \mathbin{++} bs) \, zs \, ws \\
& \equiv \{ \text{definition of Rep} \} \\
& \quad \lambda xs \rightarrow \lambda ys \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow \\
& \quad \quad \quad \mathbf{let} \, ws = (z : ys) \, \mathbf{in} \\
& \quad \quad \quad \quad (\checkmark \mathbf{Rep}[f]) \, zs \, ws \\
& \equiv \{ \text{taking this as our definition of } \mathbb{G} \} \\
& \quad \mathbb{G}[\checkmark \mathbf{Rep}[f]]
\end{aligned}$$

The step marked \star is valid because $\checkmark[-]$ is itself an evaluation context, being syntactic sugar for $\mathbf{let} \, x = [-] \, \mathbf{in} \, x$. Thus we have derived a definition of \mathbb{G} , from which we create the following factorised program:

$$\begin{aligned}
& \mathbf{reverse} = \mathbf{let} \, \{ \mathbf{rec} = \mathbb{G}[\mathbf{rec}] \} \, \mathbf{in} \, \mathbf{Abs}[\mathbf{rec}] \\
& \mathbb{G}[\mathbf{rec}] = \lambda xs \rightarrow \lambda ys \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow \mathbf{let} \, ws = (z : ys) \, \mathbf{in} \\
& \quad \quad \quad \mathbf{rec} \, zs \, ws
\end{aligned}$$

Expanding this out, we obtain:

$$\begin{aligned}
& \mathbf{reverse} = \mathbf{let} \, \{ \mathbf{rec} = \\
& \quad \lambda xs \rightarrow \lambda ys \rightarrow \mathbf{case} \, xs \, \mathbf{of} \\
& \quad \quad [] \rightarrow ys \\
& \quad \quad (z : zs) \rightarrow \mathbf{let} \, ws = (z : ys) \, \mathbf{in} \\
& \quad \quad \quad \mathbf{rec} \, zs \, ws \} \\
& \quad \mathbf{in} \, \lambda xs \rightarrow \mathbf{rec} \, xs \, []
\end{aligned}$$

The result is an implementation of fast reverse as a recursive let. The calculations here have essentially the same structure as the correctness proofs, with the addition of some administrative steps to do with the manipulation of ticks.

To illustrate the performance gain, we have graphed the performance of the original *reverse* function against the optimised version in Figure 2. We used the Criterion benchmarking library [18] with a range of list lengths to compare the performance of the two functions. The resulting graph shows a clear improvement from quadratic time to linear. We chose to use relatively small list lengths for our graphs, but the trend continues for larger values.

6.2 Tabulating a Function

Our second example is that of tabulating a function by producing a stream (infinite list) of results. Given a function f that takes a natural number as its argument, the *tabulate* function should produce the following result:

$$[f0, f1, f2, f3, \dots]$$

This function can be implemented in Haskell as follows:

$$\mathbf{tabulate} \, f = f0 : \mathbf{tabulate} \, (f \circ (+1))$$

This definition is inefficient, as it requires that the argument to f be recalculated for each element of the result stream. Essentially, this definition corresponds to the following calculation, involving a significant amount of repeated work:

$$[f0, f(0+1), f((0+1)+1), f(((0+1)+1)+1), \dots]$$

We wish to apply the worker/wrapper technique to improve the time performance of this program. The first step is to write it as a recursive let in our language:

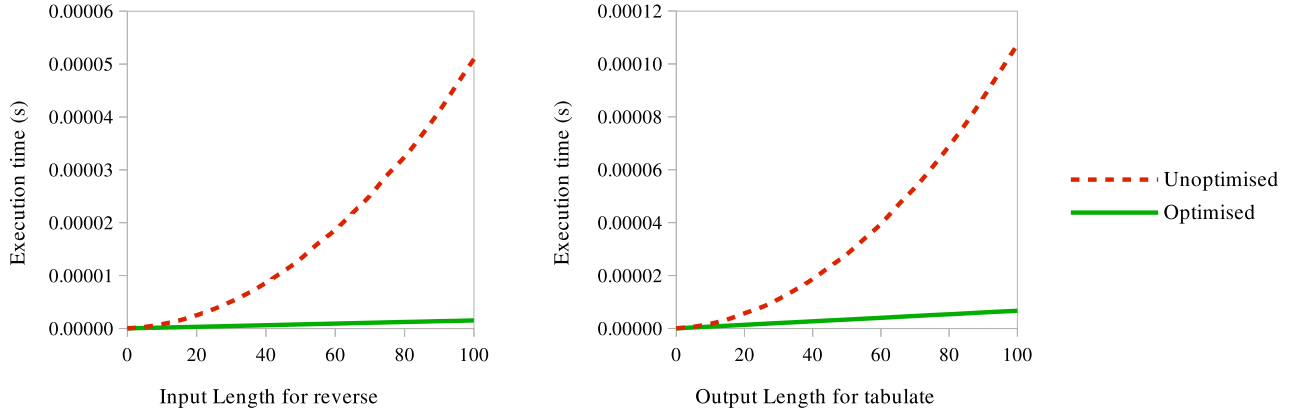


Figure 2. Performance comparisons of *reverse* and *tabulate*

$tabulate = \text{let } \{h = \mathbb{F}[h]\} \text{ in } h$
 $\mathbb{F}[M] = \lambda f \rightarrow \text{let } \{f' = \lambda x \rightarrow$
 $\quad \text{let } \{x' = x + 1\} \text{ in } f x'\}$
 $\quad \text{in } f 0 : M f'$

Next, we must devise **Abs** and **Rep** contexts. In order to avoid the repeated work, we hope to derive a version of the *tabulate* function that takes an additional number argument telling it where to “start” from. The following **Abs** and **Rep** contexts convert between these two versions:

$\text{Abs}[M] = \lambda f \rightarrow M 0 f$
 $\text{Rep}[M] = \lambda n \rightarrow \lambda f \rightarrow \text{let } \{f' = \lambda x \rightarrow$
 $\quad \text{let } \{x' = x + n\}$
 $\quad \text{in } f x'\}$
 $\quad \text{in } M f'$

Once again, we must introduce some new rules before we can derive the factorised program. Firstly, we require the following two *variable substitution* rules from [15]:

$\text{let } \{x = y\} \text{ in } \mathbb{C}[x] \triangleright \text{let } \{x = y\} \text{ in } \mathbb{C}[y]$
 $\text{let } \{x = y\} \text{ in } \mathbb{C}[y] \trianglelefteq \text{let } \{x = y\} \text{ in } \mathbb{C}[x]$

Next, we must use some properties of addition. Firstly, we have the following *identity* properties:

$x + 0 \trianglelefteq x$
 $0 + x \trianglelefteq x$

We also use the following property, combining associativity and commutativity. We shall refer to this as *associativity of +*. Where t is not free in \mathbb{C} , we have:

$\text{let } \{t = x + y\} \text{ in}$
 $\quad \text{let } \{r = t + z\} \text{ in } \mathbb{C}[r]$
 \trianglelefteq
 $\text{let } \{t = z + y\} \text{ in}$
 $\quad \text{let } \{r = x + t\} \text{ in } \mathbb{C}[r]$

Finally, we use the fact that sums may be floated out of arbitrary contexts. Where z does not occur in \mathbb{C} , we have:

$\mathbb{C}[\text{let } \{z = y + x\} \text{ in } M] \trianglelefteq \text{let } \{z = y + x\} \text{ in } \mathbb{C}[M]$

Now we can begin to apply worker/wrapper. Firstly, we verify that **Abs** and **Rep** satisfy assumption (B). Again, this is relatively straightforward:

$\text{Abs}[\text{Rep}[\mathbb{F}[h]]]$
 $\equiv \{ \text{definitions} \}$
 $\lambda f \rightarrow (\lambda n \rightarrow \lambda f \rightarrow \text{let } \{f' = \lambda x \rightarrow$
 $\quad \text{let } \{x' = x + n\}$
 $\quad \text{in } f x'\}$
 $\quad \text{in } \mathbb{F}[h] f) 0 f'$
 $\trianglelefteq \{ \beta\text{-reduction} \}$
 $\lambda f \rightarrow \text{let } \{f' = \lambda x \rightarrow$
 $\quad \text{let } \{x' = x + 0\}$
 $\quad \text{in } f x'\}$
 $\quad \text{in } \mathbb{F}[h] f'$
 $\trianglelefteq \{ x + 0 \trianglelefteq x \}$
 $\approx \lambda f \rightarrow \text{let } \{f' = \lambda x \rightarrow$
 $\quad \text{let } \{x' = x\}$
 $\quad \text{in } f x'\}$
 $\quad \text{in } \mathbb{F}[h] f'$
 $\equiv \{ \text{definition of } \mathbb{F} \}$
 $\lambda f \rightarrow \text{let } \{f' = \lambda x \rightarrow f x\}$
 $\quad \text{in } (\lambda f \rightarrow \text{let } \{f'' = \lambda x \rightarrow$
 $\quad \quad \text{let } \{x' = x + 1\} \text{ in } f x\}$
 $\quad \quad \text{in } f 0 : h f'') f'$
 $\trianglelefteq \{ \beta\text{-reduction} \}$
 $\lambda f \rightarrow \text{let } \{f' = \lambda x \rightarrow f x\}$
 $\quad \text{in } \text{let } \{f'' = \lambda x \rightarrow$
 $\quad \quad \text{let } \{x' = x + 1\} \text{ in } f' x'\}$
 $\quad \quad \text{in } f 0 : h f''\}$
 $\trianglelefteq \{ \text{value-}\beta \text{ on } f' \}$
 $\approx \lambda f \rightarrow \text{let } \{f'' = \lambda x \rightarrow$
 $\quad \text{let } \{x' = x + 1\} \text{ in } (\lambda x \rightarrow f x) x'\}$
 $\quad \text{in } (\lambda x \rightarrow f x) 0 : h f''\}$
 $\trianglelefteq \{ \beta\text{-reduction} \}$
 $\lambda f \rightarrow \text{let } \{f'' = \lambda x \rightarrow$
 $\quad \text{let } \{x' = x + 1\} \text{ in } f x'\}$
 $\quad \text{in } f 0 : h f''\}$

$$\equiv \{ \text{definition of } \mathbb{F} \}$$

$$\mathbb{F}[h]$$

Now we use condition (2) to derive the new definition of *tabulate*. This requires the use of a number of the properties that we presented earlier:

$$\begin{aligned} & \text{Rep}[\sqrt{\mathbb{F}[h]}] \\ \equiv & \{ \text{definitions} \} \\ & \lambda n \rightarrow \lambda f \rightarrow \text{let } \{ f' = \lambda x \rightarrow \\ & \quad \text{let } \{ x' = x + n \} \\ & \quad \text{in } f x' \} \\ & \text{in } (\wedge f \rightarrow \text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \text{let } \{ x'' = x + 1 \} \text{ in } f x'' \} \\ & \quad \text{in } f 0 : h f'') f' \\ \oplus & \{ \text{tick floating } [-] f' \} \\ & \lambda n \rightarrow \lambda f \rightarrow \text{let } \{ f' = \lambda x \rightarrow \\ & \quad \text{let } \{ x' = x + n \} \\ & \quad \text{in } f x' \} \\ & \text{in } (\wedge f \rightarrow \text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \text{let } \{ x'' = x + 1 \} \text{ in } f x'' \} \\ & \quad \text{in } f 0 : h f'') f' \\ \oplus & \{ \beta\text{-reduction} \} \\ & \lambda n \rightarrow \lambda f \rightarrow \text{let } \{ f' = \lambda x \rightarrow \\ & \quad \text{let } \{ x' = x + n \} \\ & \quad \text{in } f x' \} \\ & \text{in } \wedge \text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \text{let } \{ x' = x + 1 \} \text{ in } f' x'' \} \\ & \quad \text{in } f' 0 : h f'' \\ \oplus & \{ \text{value-}\beta \text{ on } f', \text{ garbage collection} \} \\ & \lambda n \rightarrow \lambda f \rightarrow \sqrt{\text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \text{let } \{ x' = x + 1 \} \text{ in } \\ & \quad (\wedge \lambda x \rightarrow \\ & \quad \quad \text{let } \{ x'' = x + n \} \\ & \quad \quad \text{in } f x'' \} x' \} \\ & \quad \text{in } (\wedge \lambda x \rightarrow \text{let } \{ x'' = x + n \} \\ & \quad \text{in } f x'') 0 : h f''} \\ \ni & \{ \text{removing ticks, } \beta\text{-reduction} \} \\ & \lambda n \rightarrow \lambda f \rightarrow \sqrt{\text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \text{let } \{ x' = x + 1 \} \text{ in } \\ & \quad \text{let } \{ x'' = x' + n \} \\ & \quad \text{in } f x'' \} \\ & \quad \text{in } (\text{let } \{ x'' = 0 + n \} \\ & \quad \text{in } f x'') : h f''} \\ \oplus & \{ \text{associativity and identity of } + \} \\ & \lambda n \rightarrow \lambda f \rightarrow \sqrt{\text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \text{let } \{ n' = n + 1 \} \text{ in } \\ & \quad \text{let } \{ x'' = x + n' \} \\ & \quad \text{in } f x'' \} \\ & \quad \text{in } (\text{let } \{ x'' = n \} \\ & \quad \text{in } f x'') : h f''} \\ \ni & \{ \text{variable substitution, garbage collection} \} \\ & \lambda n \rightarrow \lambda f \rightarrow \sqrt{\text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \text{let } \{ n' = n + 1 \} \text{ in } \\ & \quad \text{let } \{ x'' = x + n' \} \\ & \quad \text{in } f x'' \} \\ & \quad \text{in } f n : h f''} \\ \oplus & \{ \text{value let-floating} \} \\ & \lambda n \rightarrow \lambda f \rightarrow f n : \\ & \quad \sqrt{\text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \quad \text{let } \{ n' = n + 1 \} \text{ in } \\ & \quad \quad \text{let } \{ x'' = x + n' \} \\ & \quad \quad \text{in } f x'' \} \end{aligned}$$

$$\begin{aligned} & \text{in } h f'' \\ \oplus & \{ \text{sums float} \} \\ & \lambda n \rightarrow \lambda f \rightarrow f n : \\ & \quad \text{let } \{ n' = n + 1 \} \text{ in } \\ & \quad \quad \sqrt{\text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \quad \quad \text{let } \{ x'' = x + n' \} \\ & \quad \quad \text{in } f x'' \} \\ & \quad \quad \text{in } h f''} \\ \oplus & \{ \beta\text{-expansion, tick floating} \} \\ & \lambda n \rightarrow \lambda f \rightarrow f n : \\ & \quad \text{let } \{ n' = n + 1 \} \text{ in } \\ & \quad (\wedge \lambda n \rightarrow \lambda f \rightarrow \text{let } \{ f'' = \lambda x \rightarrow \\ & \quad \quad \text{let } \{ x' = x + n \} \\ & \quad \quad \text{in } f x' \} \\ & \quad \quad \text{in } h f'') n' f \\ \equiv & \{ \text{definition of Rep} \} \\ & \lambda n \rightarrow \lambda f \rightarrow f n : \\ & \quad \text{let } \{ n' = n + 1 \} \text{ in } \\ & \quad (\sqrt{\text{Rep}[h]}) n' f \\ \equiv & \{ \text{taking this as our definition of } \mathbb{G} \} \\ & \mathbb{G}[\sqrt{\text{Rep}[h]}] \end{aligned}$$

Thus we have derived a definition of \mathbb{G} , from which we create the following factorised program:

$$\begin{aligned} \text{tabulate} &= \text{let } \{ h = \mathbb{G}[h] \} \text{ in Abs}[h] \\ \mathbb{G}[M] &= \lambda n \rightarrow \lambda f \rightarrow f n : \text{let } \{ n' = n + 1 \} \text{ in } M n' f \end{aligned}$$

This is the same optimised *tabulate* function that was proved correct in [10], and the proofs here have a similar structure to the correctness proofs from that paper, except that we have now formalised that the new version of the *tabulate* function is indeed a time improvement of the original version. We note that the proof of (B) is complicated by the fact that η -reduction is not valid in this setting. In fact, if we assumed η -reduction then our proof of (B) here could be adapted into a proof of (A).

We demonstrate the performance gain in Figure 2, again based on Criterion benchmarks. This time, we keep the same input (in this case the function $\lambda n \rightarrow n * n$), but vary how many elements of the result stream we evaluate. Once again, we have an improvement from quadratic to linear performance, and the trend continues for larger values.

7. Related Work

We divide the related work into three sections. Firstly, we discuss various approaches to the operational semantics of lazy languages. Secondly, we discuss the history of improvement theory. Finally, we discuss other approaches that have been used to formally reason about efficiency.

7.1 Lazy Operational Semantics

The notion of call-by-need evaluation was first introduced in 1971 by Wadsworth [30]. However, the semantics most widely regarded as the definition of call-by-need is the natural semantics due to Launchbury [14], which was later used by Sestoft to derive the virtual machine semantics we use in this paper [27]. Ariola, Felleisen, Maraist, Odersky and Wadler presented a *call-by-need lambda calculus* [1], with operational semantics based on reductions between terms in the source language. This calculus supports an equational theory. However, Moran and Sands showed that this equational theory is subsumed by weak cost-equivalence [15].

7.2 Improvement Theory

Improvement theory was originally developed in 1991 by Sands [21], and applied in a call-by-name setting. In 1997 this was generalised to a wide class of call-by-name and call-by-value languages, also by Sands [22]. This theory was also applicable to a general class of *resources*, rather than just space and time. The theory for lazy languages was developed by Moran and Sands for time efficiency [15] and Gustavsson and Sands for space efficiency [8, 9]. Since the last of these papers was published in 2001, there does not seem to have been much work on improvement theory. We hope that this paper can help to regenerate interest in this topic.

7.3 Formal Reasoning About Efficiency

Okasaki [17] uses techniques of amortised cost analysis to reason about the asymptotic time complexity of lazy functional data structures. This is achieved by modifying analysis techniques such as the Banker’s Method, where the notion of *credit* is used to spread out the notional cost of an expensive but infrequent operations over more frequent and cheaper operations. The key idea in Okasaki’s work is to invert such techniques to use the notion of *debt*. This allows the analyses to deal with the persistence of data structures, where the same structure may exist in multiple versions at once. While credit may only be spent once, a single debt may be paid off multiple times (in different versions of the same structure) without risking bankruptcy. These techniques have been used to analyse the asymptotic performance of a number of functional data structures.

Sansom and Peyton Jones [24] give a presentation of the GHC profiler, which can be used to measure time as well as space usage of Haskell programs. In doing so, they give a formal cost semantics for GHC Core programs based around the notion of *cost centres*. Cost centres are a way of annotating expressions, so that the profiler can indicate which parts of the source program cost the most to execute. The cost semantics is used as a specification to develop a precise profiling framework, as well as to prove various properties about cost attribution and verify that certain program transformations do not affect the attribution of costs, though they may of course reduce cost overall. Cost centres are now widely-used in profiling Haskell programs.

Hope [11] applies a technique based on *instrumenting* an abstract machine with cost information to derive a cost semantics for call-by-value functional programs. More specifically, starting from a denotational semantics for the source language, one derives an abstract machine for this language using standard program transformation techniques, instruments this machine with cost information, and then reverses the derivation to arrive at an instrumented denotational semantics. This semantics can then be used to reason about the cost of programs in the high-level source language without reference to the details of the abstract machine. This approach was used to calculate the space and time cost of a range of programming examples, as well as to derive a new deforestation theorem for hylomorphisms.

8. Conclusion

In this paper, we have shown how improvement theory can be used to justify the worker/wrapper transformation as a program optimisation, by formally proving that, under certain natural conditions, the transformation is guaranteed to preserve or improve time performance. This guarantee is with respect to an established operational semantics for

call-by-need evaluation. We then verified that two examples from previous worker/wrapper papers met the preconditions for this performance guarantee, demonstrating the use of our theory while also verifying the validity of the examples. This work appears to be the first time that rigorous performance guarantees have been given for a general purpose optimisation technique in a call-by-need setting.

8.1 Further Work

As well as for fixed points, worker/wrapper theories also exist for more structured recursion operators such as folds [13] and unfolds [10]. Though the theory we present here can be specialised to such operators, it may be beneficial to investigate this more closely, as doing so may reveal more interesting and subtle details yet to be uncovered.

As we mentioned earlier in this paper, a typed theory would be more useful, allowing more power when reasoning about programs. This would also match more closely with the original worker/wrapper theories, which were typed. The key barrier to this is that there is currently no typed improvement theory, so such a theory would have to be developed before the theory here could be made typed.

The theory we present here only applies to time efficiency. Gustavsson and Sands have developed an improvement theory for space [8, 9], so this would be an obvious next step for developing our theory. More generally, we could apply a technique such as that used by Sands [22] to develop a theory that applies to a large class of resources, and examine which assumptions must be made about the resources we consider for our theory to apply.

Assumptions (A), (B) and (C) are written as weak cost-equivalences, which limits the scope of our theory to cases where **Abs** and **Rep** are fairly simple. We would like to also be able to cover cases where the **Abs** and **Rep** contexts correspond to expensive operations, but the extra cost is made up for by the overall efficiency gain of the transformation. To cover such cases, we would require a richer version of improvement theory that is able to quantify *how much* better one program is than another.

As our examples show, the calculations required to derive an improved program can often be quite involved. The HERMIT system, devised by a team at the University of Kansas [6, 26], facilitates program transformations by providing an interactive interface for program transformation that verifies correctness. If improvement theory could be integrated into such a system, it would be significantly easier to apply our worker/wrapper improvement theory.

Finally, we are working on a general worker/wrapper theory that will apply to any operator with the property of *dinaturality* [5]. It is also interesting to consider whether such a general categorical approach can be applied to an operational theory. If this is the case, dinaturality may also provide the necessary machinery to *unify* the denotational (correctness) and operational (efficiency) theories, which as we have already observed in this paper are very similar in terms of their formulations and proofs. Voigtländer and Johann used parametricity to justify program transformations from a perspective of *observational approximation* [29]. It may be productive to investigate whether their techniques can be applied to a notion of improvement.

Acknowledgments

The authors would like to thank the reviewers as well as Neil Sculthorpe, for their helpful comments on this paper.

References

- [1] Z. M. Ariola, M. Felleisen, J. Maraist, M. Odersky, and P. Wadler. The Call-by-Need Lambda Calculus. In *POPL '95*, pages 233–246. ACM, 1995.
- [2] R. M. Burstall and J. Darlington. A Transformation System for Developing Recursive Programs. *Journal of the ACM*, 24(1):44–67, 1977.
- [3] K. Claessen and J. Hughes. QuickCheck: A Lightweight Tool for Random Testing of Haskell Programs. In *ICFP*. ACM, 2000.
- [4] T. Coquand. Infinite Objects in Type Theory. In *TYPES '93*, volume 806 of *Lecture Notes in Computer Science*. Springer, 1993.
- [5] E. Dubuc and R. Street. Dinatural Transformations. In *Lecture Notes in Mathematics*, volume 137 of *Lecture Notes in Mathematics*, pages 126–137. Springer Berlin Heidelberg, 1970.
- [6] A. Farmer, A. Gill, E. Komp, and N. Sculthorpe. The HERMIT in the Machine: A Plugin for the Interactive Transformation of GHC Core Language Programs. In *Haskell Symposium (Haskell '12)*, pages 1–12. ACM, 2012.
- [7] A. Gill and G. Hutton. The Worker/Wrapper Transformation. *Journal of Functional Programming*, 19(2), 2009.
- [8] J. Gustavsson and D. Sands. A Foundation for Space-Safe Transformations of Call-by-Need Programs. *Electronic Notes on Theoretical Computer Science*, 26:69–86, 1999.
- [9] J. Gustavsson and D. Sands. Possibilities and Limitations of Call-by-Need Space Improvement. In *ICFP*, pages 265–276. ACM, 2001.
- [10] J. Hackett, G. Hutton, and M. Jaskelioff. The Under Performing Unfold: A New Approach to Optimising Corecursive Programs. 2013. To appear in the volume of selected papers from the 25th International Symposium on Implementation and Application of Functional Languages, Nijmegen, The Netherlands, August 2013.
- [11] C. Hope. *A Functional Semantics for Space and Time*. PhD thesis, School of Computer Science, University of Nottingham, 2008.
- [12] J. Hughes. A Novel Representation of Lists and its Application to the Function “reverse”. *Information Processing Letters*, 22(3):141–144, 1986.
- [13] G. Hutton, M. Jaskelioff, and A. Gill. Factorising Folds for Faster Functions. *Journal of Functional Programming Special Issue on Generic Programming*, 20(3&4), June 2010.
- [14] J. Launchbury. A Natural Semantics for Lazy Evaluation. In *POPL '93*, pages 144–154. ACM, 1993.
- [15] A. Moran and D. Sands. Improvement in a Lazy Context: An Operational Theory for Call-by-Need. Extended version of [16], available at <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.37.144&rep=rep1&type=pdf>.
- [16] A. Moran and D. Sands. Improvement in a Lazy Context: An Operational Theory for Call-by-Need. In *POPL '99*, pages 43–56. ACM, 1999.
- [17] C. Okasaki. *Purely Functional Data Structures*. Cambridge University Press, 1999.
- [18] B. O’Sullivan. Criterion, A New Benchmarking Library for Haskell, 2009. URL <http://www.serpentine.com/blog/2009/09/29/criterion-a-new-benchmarking-library-for-haskell/>.
- [19] W. Partain. The nofib Benchmark Suite of Haskell Programs. In *Glasgow Workshop on Functional Programming*, pages 195–202. Springer, 1992.
- [20] S. L. Peyton Jones. Compiling Haskell by Program Transformation: A Report from the Trenches. In *ESOP*, pages 18–44. Springer, 1996.
- [21] D. Sands. Operational Theories of Improvement in Functional Languages (Extended Abstract). In *Glasgow Workshop on Functional Programming*, 1991.
- [22] D. Sands. From SOS Rules to Proof Principles: An Operational Metatheory for Functional Languages. In *POPL '97*, pages 428–441. ACM Press, 1997.
- [23] D. Sands. Improvement Theory and its Applications. In *Higher Order Operational Techniques in Semantics, Publications of the Newton Institute*. Cambridge University Press, 1997.
- [24] P. M. Sansom and S. L. Peyton Jones. Formally Based Profiling for Higher-Order Functional Languages. *ACM Transactions on Programming Languages and Systems*, 19(2), 1997.
- [25] N. Sculthorpe and G. Hutton. Work It, Wrap It, Fix It, Fold It. *Journal of Functional Programming*, 2014.
- [26] N. Sculthorpe, A. Farmer, and A. Gill. The HERMIT in the Tree: Mechanizing Program Transformations in the GHC Core Language. In *Proceedings of Implementation and Application of Functional Languages (IFL '12)*, volume 8241 of *Lecture Notes in Computer Science*, pages 86–103, 2013.
- [27] P. Sestoft. Deriving a Lazy Abstract Machine. *Journal of Functional Programming*, 7(3):231–264, 1997.
- [28] D. A. Turner. Elementary Strong Functional Programming. In *FPLE '95*, volume 1022 of *Lecture Notes in Computer Science*. Springer, 1995.
- [29] J. Voigtländer and P. Johann. Selective Strictness and Parametricity in Structural Operational Semantics, Inequationally. *Theor. Comput. Sci.*, 388(1-3):290–318, 2007.
- [30] C. P. Wadsworth. *Semantics and Pragmatics of the Lambda Calculus*. PhD thesis, Computing Laboratory, University of Oxford, 1971.
- [31] O. Wilde. Lady Windermere’s Fan, A Play About a Good Woman. First performed in 1892.